

Original pca Slides

Dimensionality Reduction Principal Component Analysis (PCA)

CS229: Machine Learning

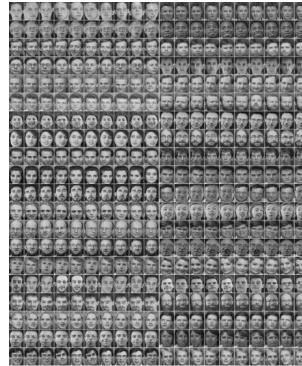
Carlos Guestrin

Stanford University

Slides include content developed by and co-developed with Emily Fox

Embedding

Example: Embedding images to visualize data



Images with
thousands or
millions of pixels

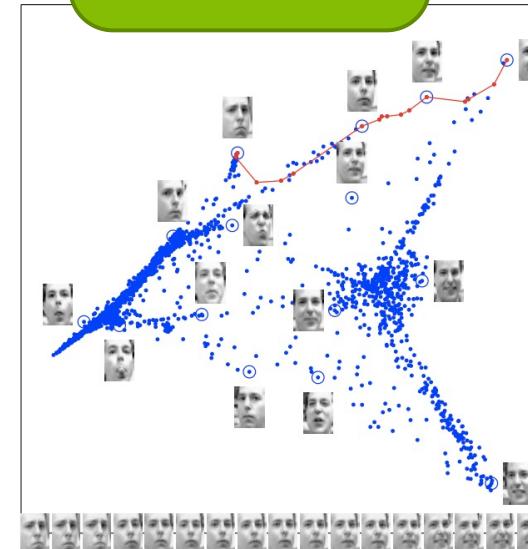


PCA

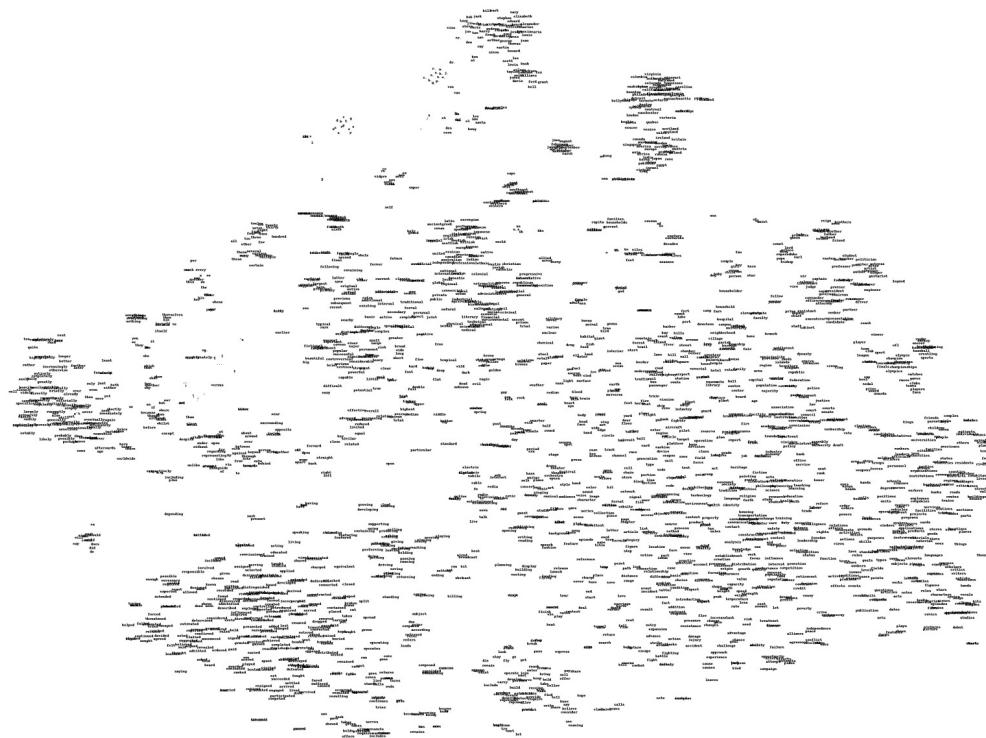


Intelligence

Can we give each
image a coordinate,
such that similar
images are near
each other?

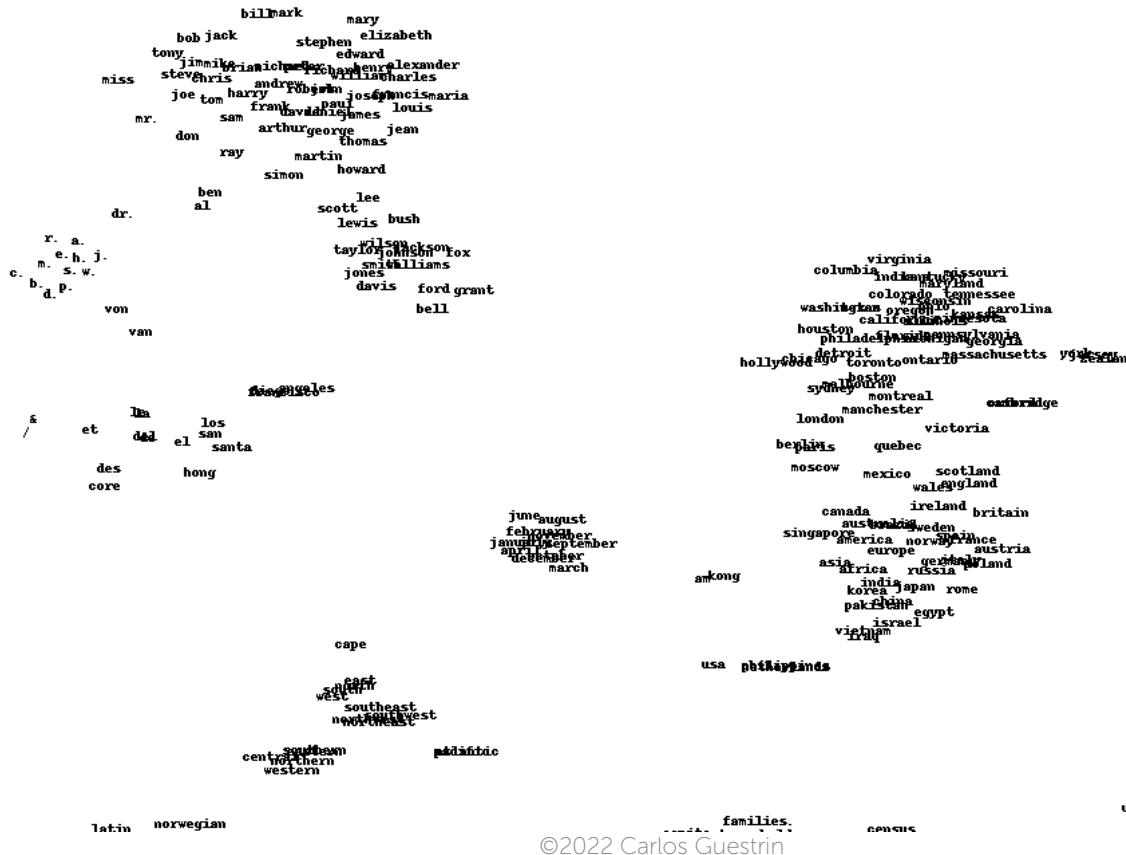


Embedding words



[Joseph Turian 2008]

Embedding words (zoom in)



[Joseph Turian 2008]

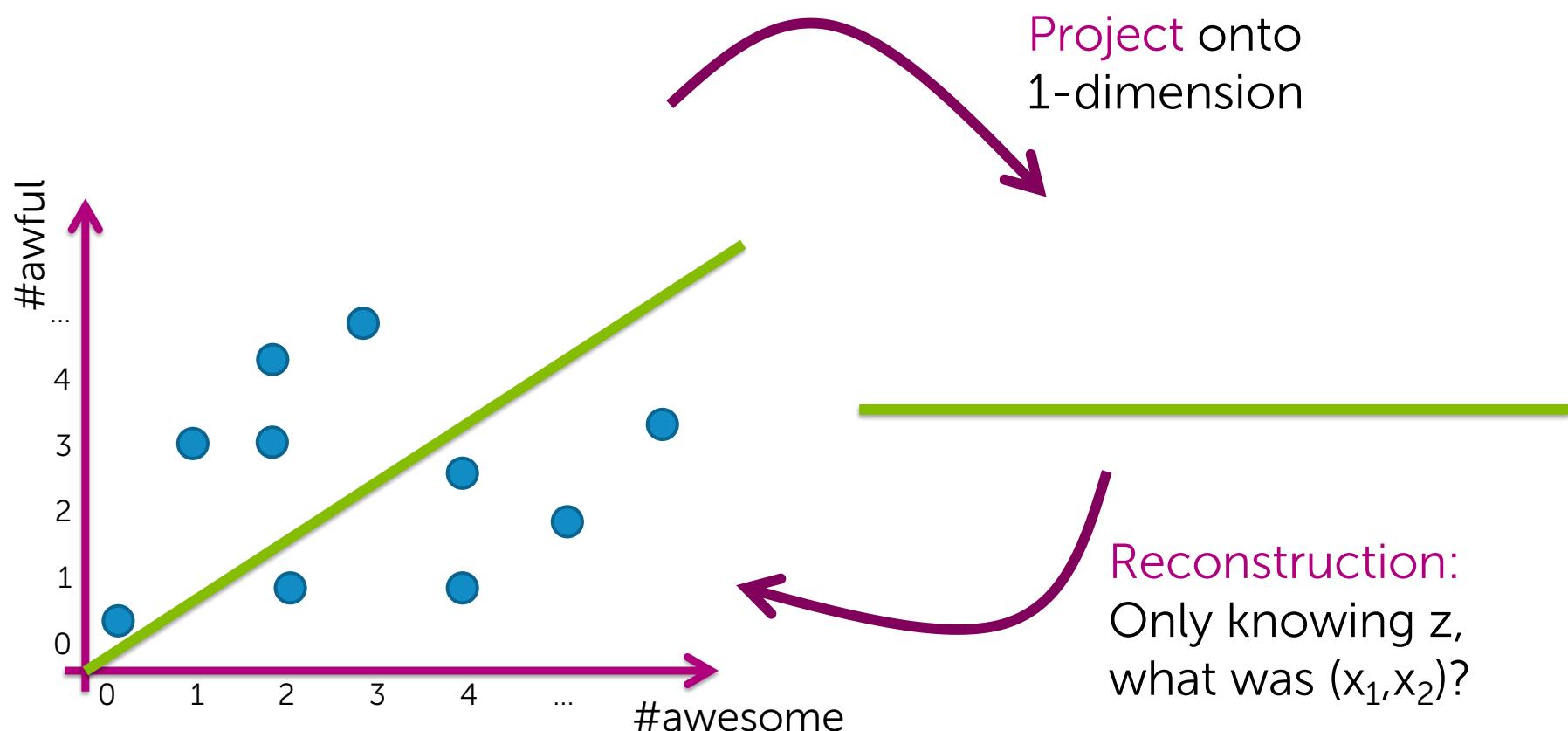
Dimensionality reduction

- Input data may have thousands or millions of dimensions!
 - e.g., text data
- **Dimensionality reduction:** represent data with fewer dimensions
 - easier learning – fewer parameters
 - visualization – hard to visualize more than 3D or 4D
 - discover “**intrinsic dimensionality**” of data
 - high dimensional data that is truly lower dimensional

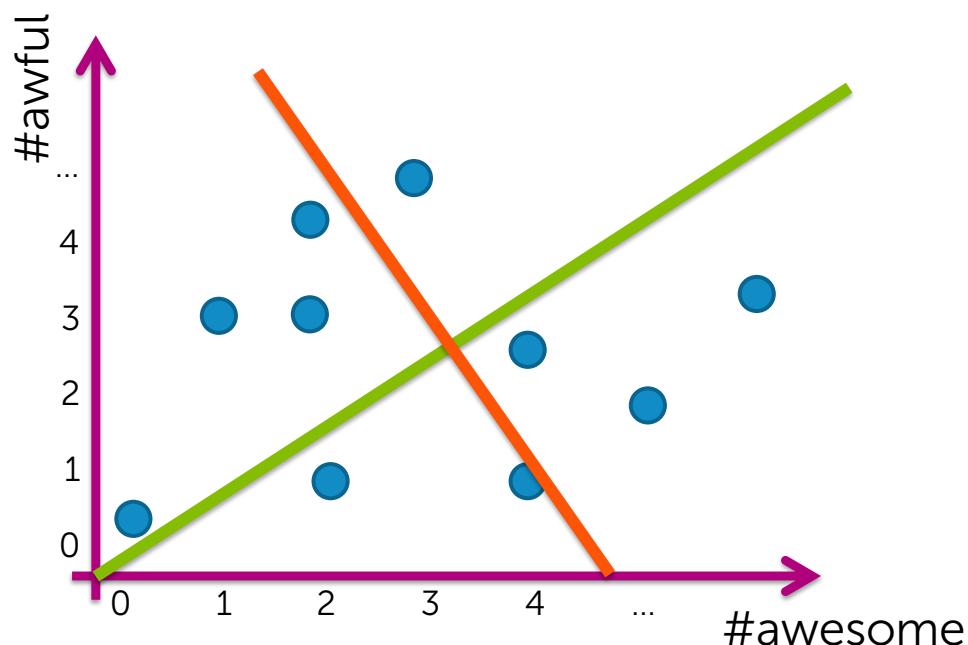
Lower dimensional projections

- Rather than picking a subset of the features, we can **create new features** that are **combinations of existing features**
- Let's see this in the unsupervised setting
 - just x , but no y

Linear projection and reconstruction



What if we project onto d vectors?



Perfect
reconstruction!

If I had to choose one of these vectors, which do I prefer?



Principal component analysis (PCA) – Basic idea

- Project d -dimensional data into k -dimensional space while preserving as much information as possible:
 - e.g., project space of 10000 words into 3-dimensions
 - e.g., project 3-d into 2-d
- Choose projection with **minimum reconstruction error**

“PCA explained visually”

<http://setosa.io/ev/principal-component-analysis/>

Linear projections, a review

- Project a point into a (lower dimensional) space:
 - point: $x = (x_1, \dots, x_d)$
 - select a basis – set of basis vectors – (u_1, \dots, u_k)
 - we consider orthonormal basis:
 - $u_i \bullet u_i = 1$, and $u_i \bullet u_j = 0$ for $i \neq j$
 - select a center – \bar{x} , defines offset of space
 - best coordinates in lower dimensional space defined by dot-products: (z_1, \dots, z_k) , $z_i = (\bar{x} - x) \bullet u_i$
 - minimum squared error
-

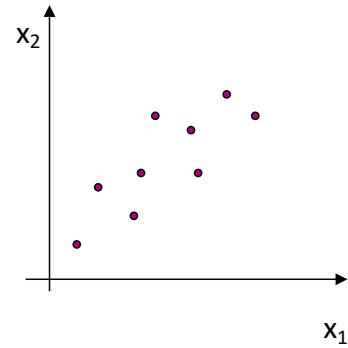
PCA finds projection that minimizes reconstruction error

- Given N data points: $\mathbf{x}^i = (x_1^i, \dots, x_d^i)$, $i=1\dots N$
- Will represent each point as a projection:

$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j \quad \text{and} \quad \bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^i \quad z_j^i = (\mathbf{x}^i - \bar{\mathbf{x}}) \cdot \mathbf{u}_j$$

- PCA:
 - Given $k \ll d$, find $(\mathbf{u}_1, \dots, \mathbf{u}_k)$ minimizing reconstruction error:

$$error_k = \sum_{i=1}^N (\mathbf{x}^i - \hat{\mathbf{x}}^i)^2$$



Understanding the reconstruction error

- Note that \mathbf{x}^i can be represented exactly by d -dimensional projection:

$$\mathbf{x}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$

- Rewriting error:

$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$

$$z_j^i = (\mathbf{x}^i - \bar{\mathbf{x}}) \cdot \mathbf{u}_j$$

Given $k \ll d$, find $(\mathbf{u}_1, \dots, \mathbf{u}_k)$
minimizing reconstruction error:

$$error_k = \sum_{i=1}^N (\mathbf{x}^i - \hat{\mathbf{x}}^i)^2$$

Reconstruction error and covariance matrix

$$error_k = \sum_{i=1}^N \sum_{j=k+1}^d [\mathbf{u}_j \cdot (\mathbf{x}^i - \bar{\mathbf{x}})]^2$$

$$\Sigma = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}^i - \bar{\mathbf{x}})(\mathbf{x}^i - \bar{\mathbf{x}})^T$$

Minimizing reconstruction error and eigen vectors

- Minimizing reconstruction error equivalent to picking orthonormal basis $(\mathbf{u}_1, \dots, \mathbf{u}_d)$ minimizing:

$$\text{error}_k = \frac{1}{N} \sum_{j=k+1}^d \mathbf{u}_j^T \Sigma \mathbf{u}_j$$

- Eigen vector:

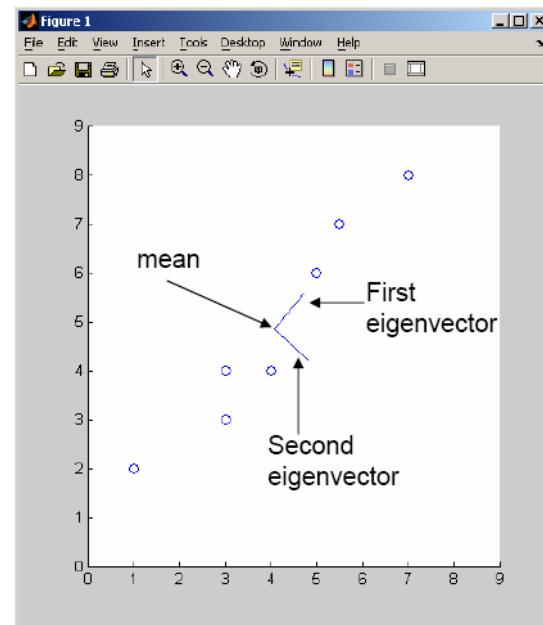
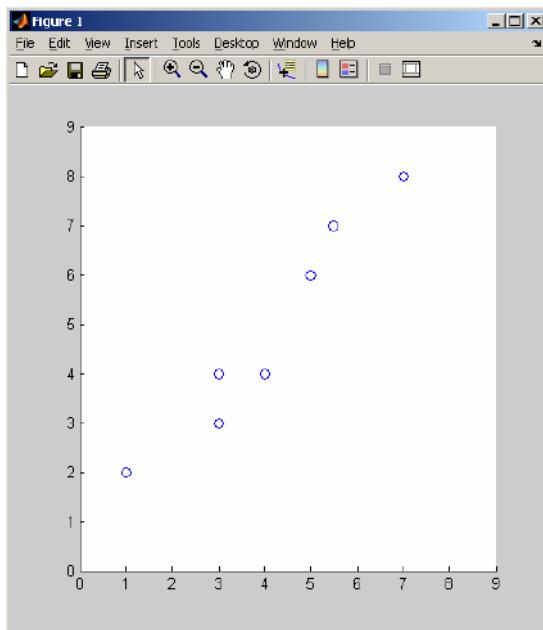
- Minimizing reconstruction error equivalent to picking $(\mathbf{u}_{k+1}, \dots, \mathbf{u}_d)$ to be eigen vectors with smallest eigen values

Basic PCA algorithm

- Start from N by d data matrix X
- **Recenter:** subtract mean from each row of X
 - $X_c \leftarrow X - \bar{X}$
- **Compute covariance matrix:**
 - $\Sigma \leftarrow 1/N X_c^T X_c$
- Find **eigen vectors and values** of Σ
- **Principal components:** k eigen vectors with highest eigen values

PCA example

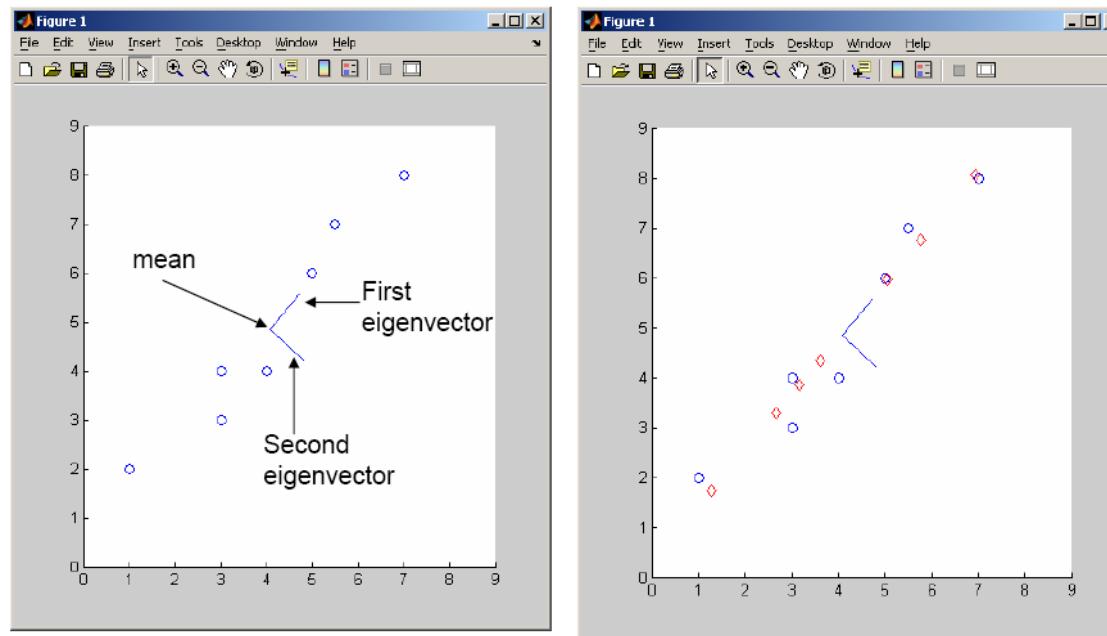
$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$



PCA example – reconstruction

$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$

only used first principal component



Eigenfaces [Turk, Pentland '91]

- Input



- Principal components:



Eigenfaces reconstruction

- Each image corresponds to adding 8 principal components:



Scaling up

- Covariance matrix can be really big!
 - Σ is d by d
 - Say, only 10000 features
 - finding eigenvectors is very slow...
- Use singular value decomposition (SVD)
 - finds k eigenvectors
 - great implementations available, e.g., python, R, Matlab svd

SVD

- Write $X = W S V^T$
 - $X \leftarrow$ data matrix, one row per datapoint
 - $W \leftarrow$ weight matrix, one row per datapoint – coordinate of x^i in eigenspace
 - $S \leftarrow$ singular value matrix, diagonal matrix
 - in our setting each entry is eigenvalue λ_j
 - $V^T \leftarrow$ singular vector matrix
 - in our setting each row is eigenvector v_j

PCA using SVD algorithm

- Start from m by n data matrix X
- **Recenter:** subtract mean from each row of X
 - $X_c \leftarrow X - \bar{X}$
- Call SVD algorithm on X_c – ask for k singular vectors
- **Principal components:** k singular vectors with highest singular values (rows of V^T)
 - **Coefficients** become:

What you need to know

- Dimensionality reduction
 - why and when it's important
- Simple feature selection
- Principal component analysis
 - minimizing reconstruction error
 - relationship to covariance matrix and eigenvectors
 - using SVD

Annotated pca Slides

Dimensionality Reduction Principal Component Analysis (PCA)

CS229: Machine Learning

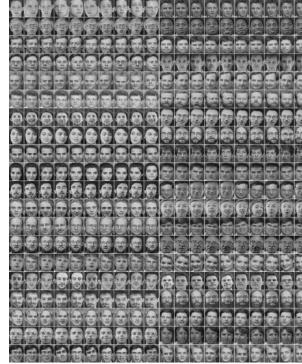
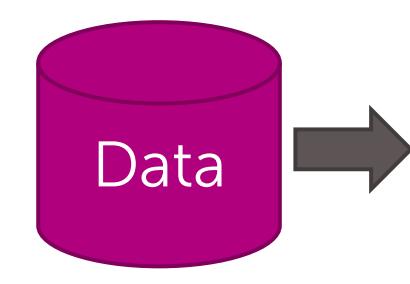
Carlos Guestrin

Stanford University

Slides include content developed by and co-developed with Emily Fox

Embedding

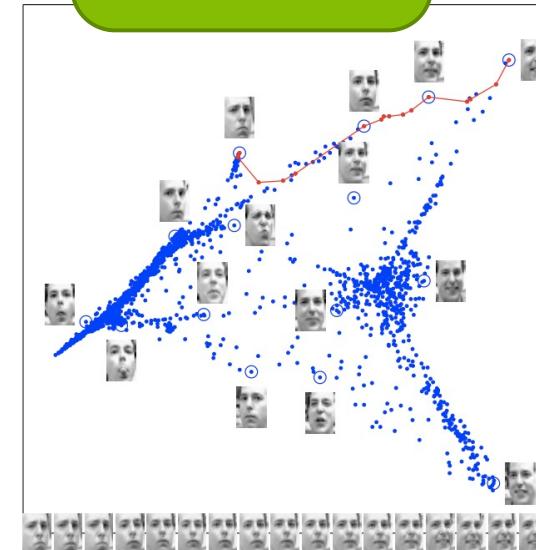
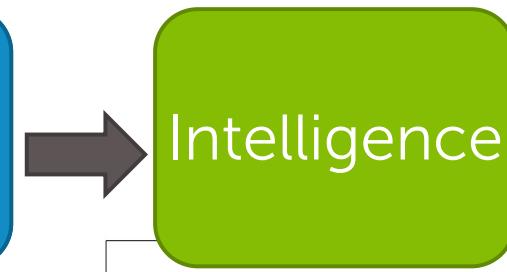
Example: Embedding images to visualize data



Images with
thousands or
millions of pixels



Can we give each
image a coordinate,
such that similar
images are near
each other?

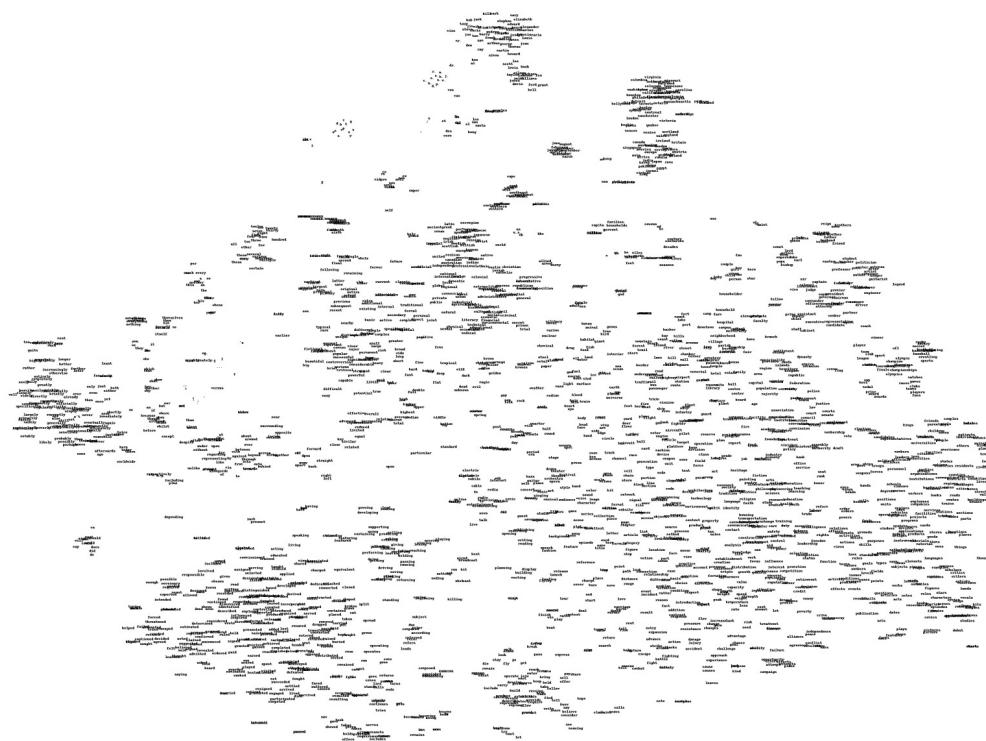


[Saul &
Roweis '03]

CS229: Machine Learning

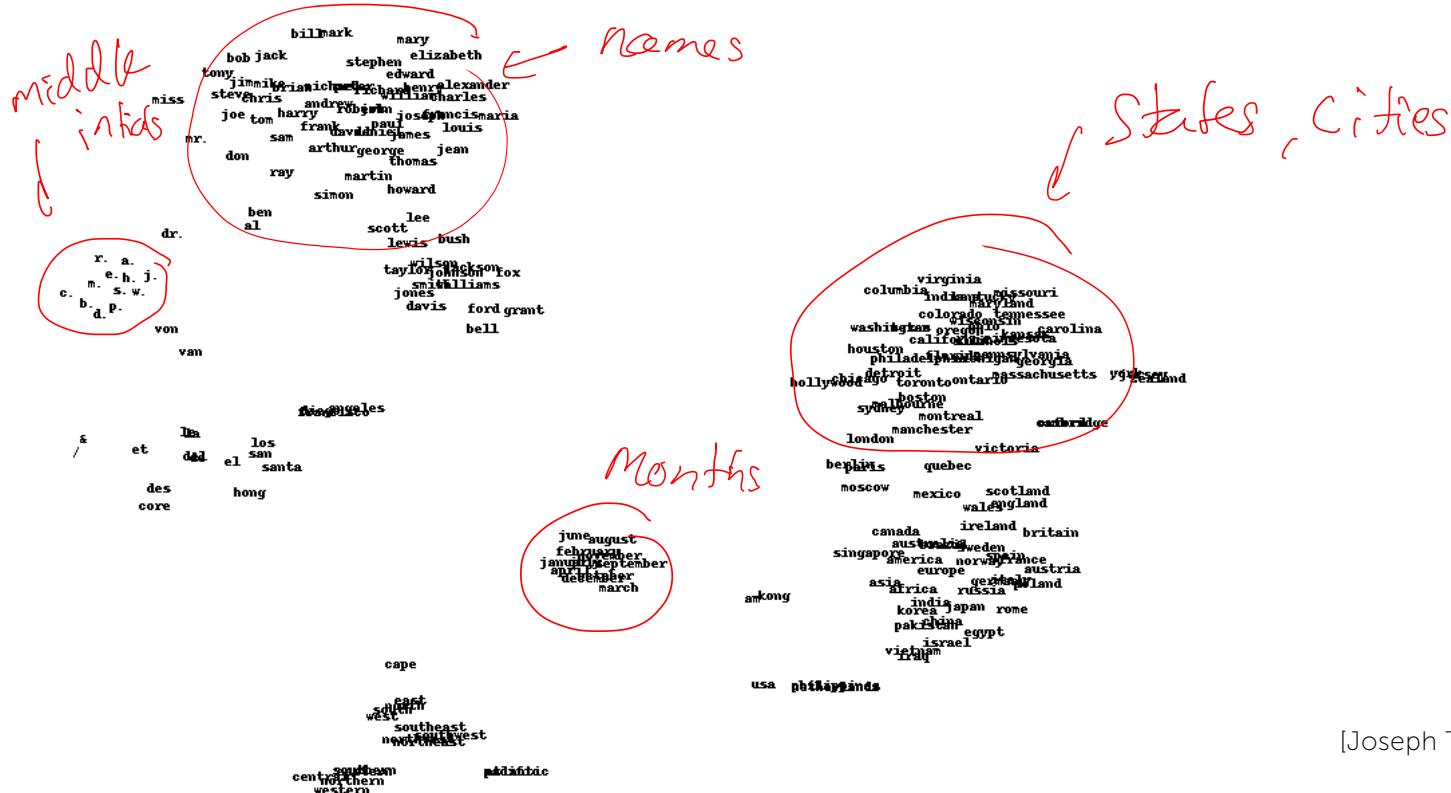
Embedding words

10,000 word dictionary



[Joseph Turian 2008]

Embedding words (zoom in)



[Joseph Turian 2008]

Dimensionality reduction

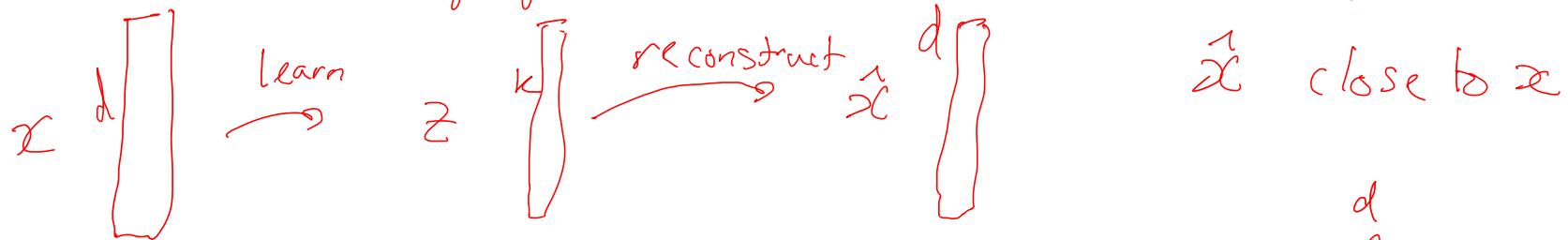
- Input data may have thousands or millions of dimensions!
 - e.g., text data
- **Dimensionality reduction:** represent data with fewer dimensions
 - easier learning – fewer parameters
 - visualization – hard to visualize more than 3D or 4D
 - discover “**intrinsic dimensionality**” of data
 - high dimensional data that is truly lower dimensional

Lower dimensional projections

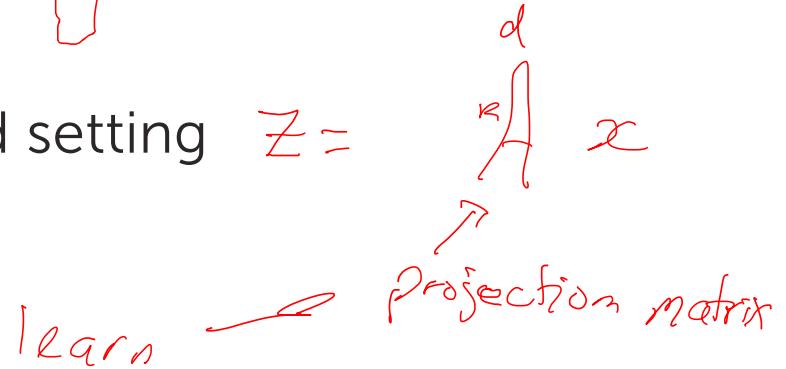
$$d \gg k$$

- Rather than picking a subset of the features, we can **create new features** that are **combinations of existing features**

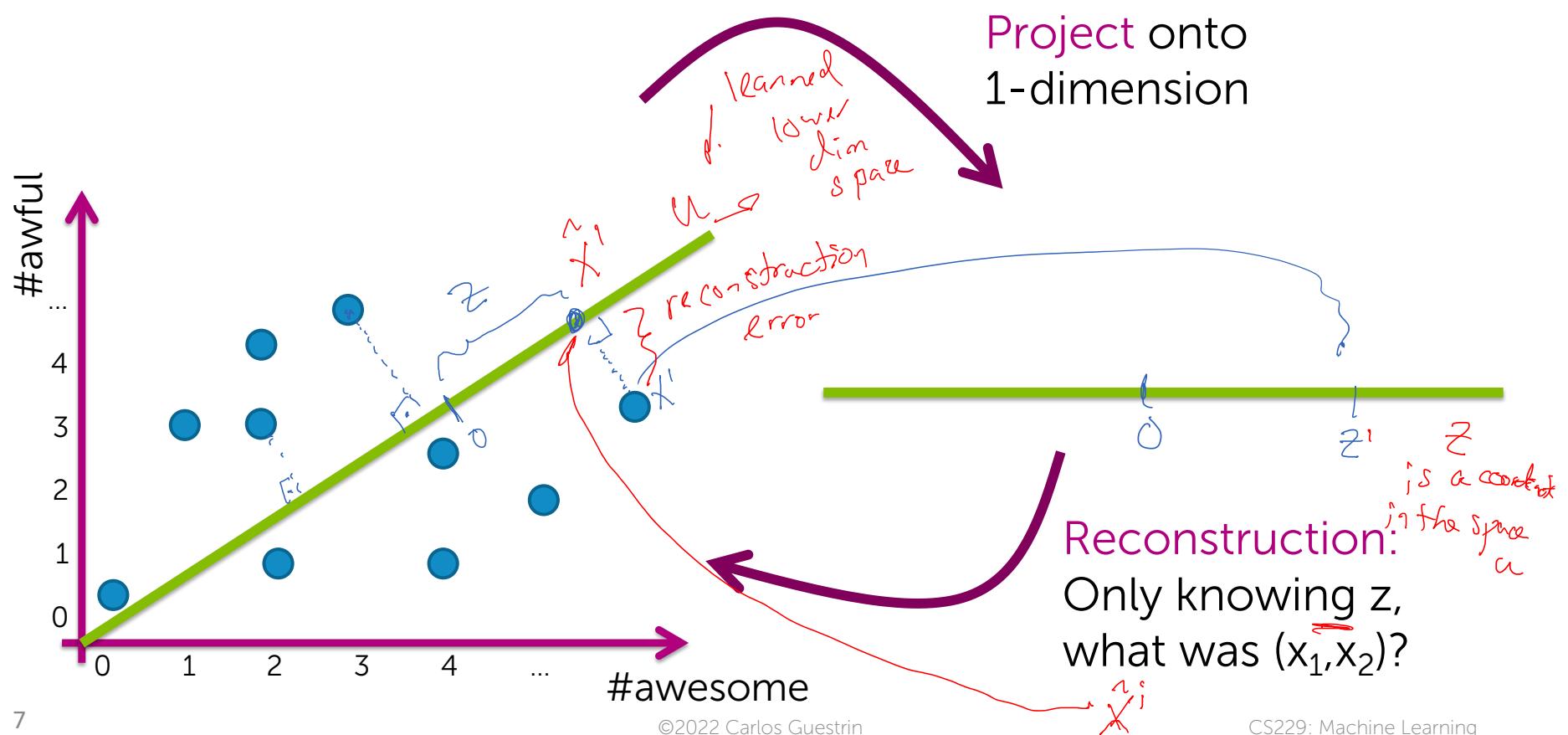
PCA: linear projections: $z_1 = 2.5x_1 - 3.2x_2 + 3.7x_3, \dots$



- Let's see this in the unsupervised setting
– just x , but no y

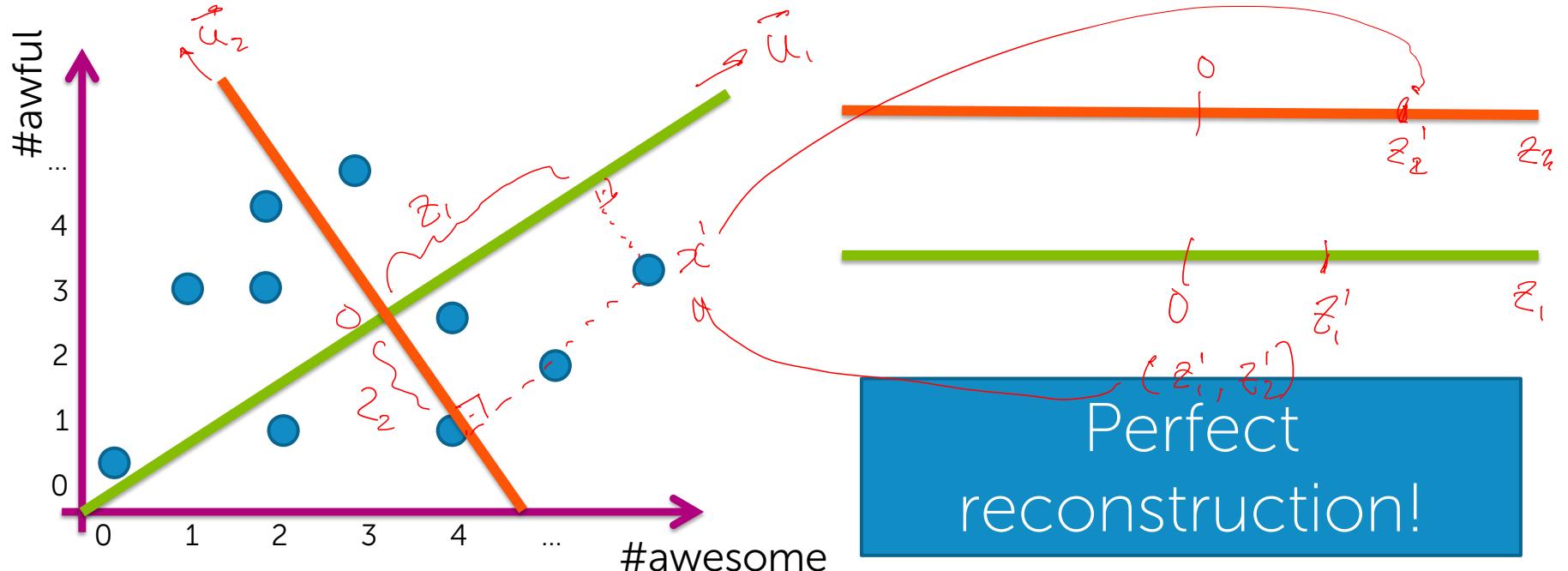


Linear projection and reconstruction

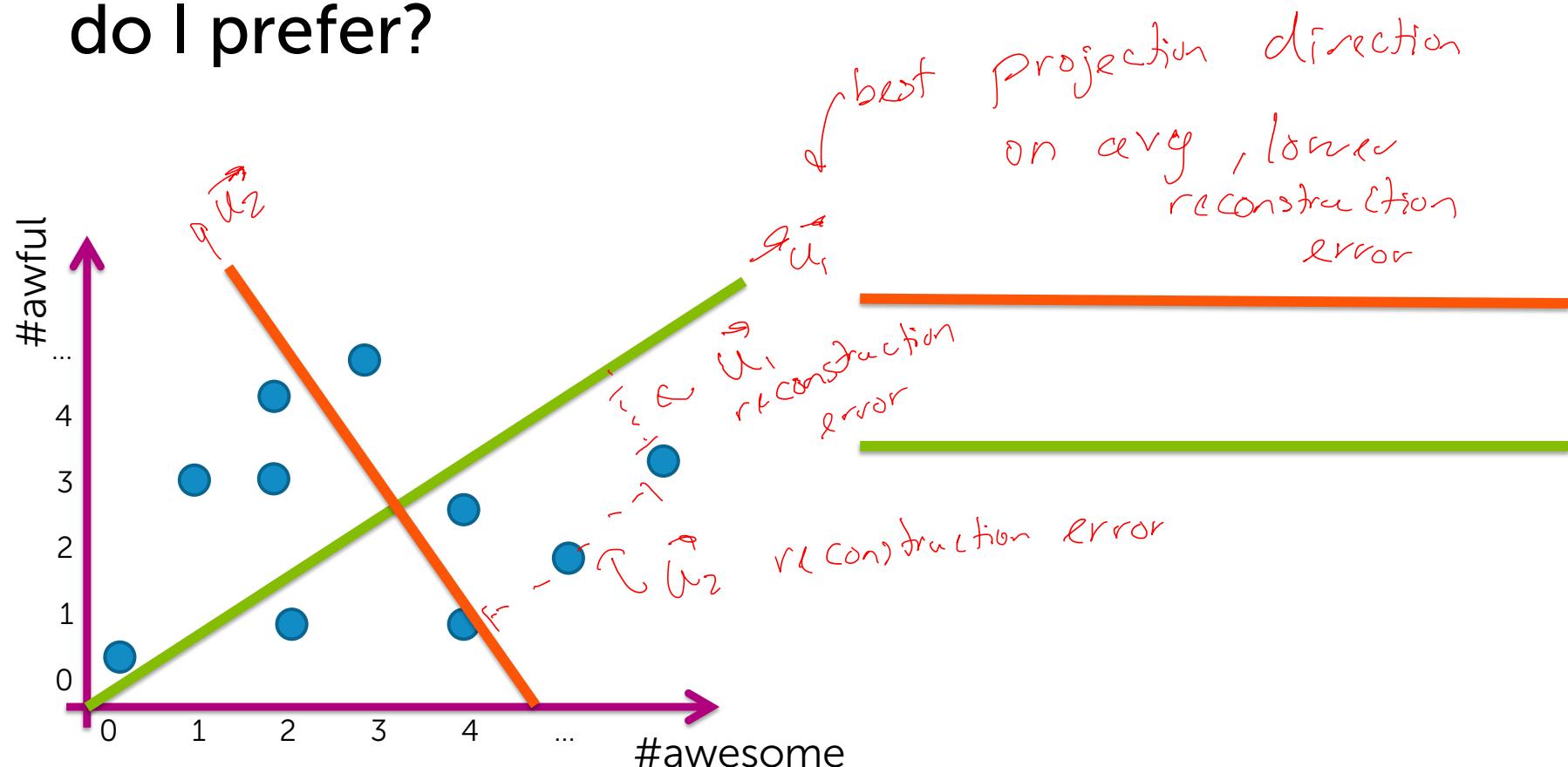


What if we project onto d vectors?

$$x' = z_1 \vec{u}_1 + z_2 \vec{u}_2 \quad (\text{ignoring offset})$$



If I had to choose one of these vectors, which do I prefer?



Principal component analysis (PCA) – Basic idea

- Project d -dimensional data into k -dimensional space while preserving as much information as possible:
 - e.g., project space of 10000 words into 3-dimensions
 - e.g., project 3-d into 2-d
- Choose projection with minimum reconstruction error

“PCA explained visually”

<http://setosa.io/ev/principal-component-analysis/>

Linear projections, a review

- Project a point into a (lower dimensional) space:
 - point: $x = (x_1, \dots, x_d)$ 10,000 dims
 - select a basis – set of basis vectors – (u_1, \dots, u_k)
 - we consider orthonormal basis: $u_i \perp u_j$
 - $u_i \cdot u_i = 1$, and $u_i \cdot u_j = 0$ for $i \neq j$
 - select a center – \bar{x} , defines offset of space mean value of x
 - best coordinates in lower dimensional space defined by dot-products:
- (z_1, \dots, z_k) , $z_i = (\bar{x} - x) \cdot u_i$
- dot product minimizes the squared error of projection
-
- \bar{x} \tilde{x} $x - \bar{x}$ x
- $$z_i = (x - \bar{x}) \cdot \vec{u}_i \quad \boxed{z_i = \underset{\bar{z}}{\operatorname{arg\min}} ((x - \bar{z}) \cdot \vec{u}_i)^2}$$

project x into u
coordinates z_1, \dots, z_k

$$\tilde{x} = \bar{x} + \sum_{j=1}^k z_j u_j$$

reconstruction

if $k = d$

$$\tilde{x} = x$$

PCA finds projection that minimizes reconstruction error

- Given N data points: $\mathbf{x}^i = (x_1^i, \dots, x_d^i)$, $i=1 \dots N$
- Will represent each point as a projection:

$$\text{Reconstructed } \hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$

and $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^i$

mean vector
avg x over data
Projection direction
coordinates

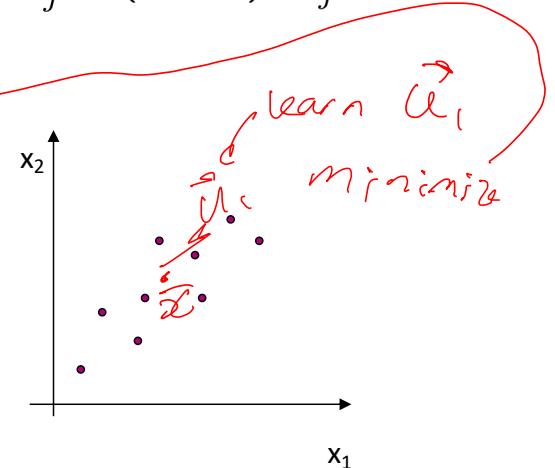
- PCA:
 - Given $k \ll d$, find $(\mathbf{u}_1, \dots, \mathbf{u}_k)$ minimizing reconstruction error:

Want minimize $\text{error}_k = \sum_{i=1}^N (\mathbf{x}^i - \hat{\mathbf{x}}^i)^2$

over choice of $\bar{\mathbf{u}}, \dots, \mathbf{u}_k$

from previous slide:
 projection coordinate is
 dot product

$$z_j^i = (\mathbf{x}^i - \bar{\mathbf{x}}) \cdot \mathbf{u}_j$$



Understanding the reconstruction error

- Note that x^i can be represented exactly by d -dimensional projection:

$$x^i = \bar{x} + \sum_{j=1}^k z_j^i u_j$$

- Rewriting error:

$$\begin{aligned} \min_{\mathbf{u}} \text{error}_k &= \min_{\mathbf{u}} \sum_{i=1}^N (x^i - \hat{x}^i)^2 = \sum_{i=1}^N \left[\bar{x} + \sum_{j=1}^d z_j^i u_j - \left(\bar{x} + \sum_{j=1}^k z_j^i u_j \right) \right]^2 \\ &= \sum_{i=1}^N \left[\sum_{j=k+1}^d z_j^i u_j \right]^2 = \sum_{i=1}^N \left[\sum_{j=k+1}^d z_j^i u_j \cdot u_j z_j^i \right] + \sum_{j=k+1}^d \sum_{j' \neq j} z_j^i u_j \cdot u_{j'} z_{j'}^i \quad \text{orthonormal} \\ &= \sum_{i=1}^N \sum_{j=k+1}^d (z_j^i)^2 \leftarrow \begin{array}{l} \text{minimizing reconstruction error} \\ \equiv \text{min squared thrown out coefficients} \end{array} \end{aligned}$$

$$\hat{x}^i = \bar{x} + \sum_{j=1}^k z_j^i u_j$$

$$z_j^i = (x^i - \bar{x}) \cdot u_j$$

Given $k \ll d$, find (u_1, \dots, u_k)
minimizing reconstruction error:

$$\text{error}_k = \sum_{i=1}^N (x^i - \hat{x}^i)^2$$

Reconstruction error and covariance matrix

$$\begin{aligned}
 & \text{error}_k = \sum_{i=1}^N \sum_{j=k+1}^d [u_j \cdot (\mathbf{x}^i - \bar{\mathbf{x}})]^2 \\
 &= \sum_{i=1}^N \sum_{j=k+1}^d u_j^\top (\mathbf{x}^i - \bar{\mathbf{x}}) (\mathbf{x}^i - \bar{\mathbf{x}})^\top u_j \\
 &\quad \text{push sum over data in} \\
 &= \sum_{j=k+1}^d u_j^\top \left[\sum_{i=1}^N (\mathbf{x}^i - \bar{\mathbf{x}}) (\mathbf{x}^i - \bar{\mathbf{x}})^\top \right] u_j \\
 &= N \sum_{j=k+1}^d u_j^\top \Sigma u_j
 \end{aligned}$$

choose
 u_j to
 minimize
 this error

matrix
notation

$$\begin{aligned}
 \Sigma &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}^i - \bar{\mathbf{x}})(\mathbf{x}^i - \bar{\mathbf{x}})^T \\
 \Sigma &= \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots \\ \sigma_{21} & \sigma_2^2 & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}
 \end{aligned}$$

Minimizing reconstruction error and eigen vectors

- Minimizing reconstruction error equivalent to picking orthonormal basis $(\mathbf{u}_1, \dots, \mathbf{u}_d)$ minimizing $\text{error}_k = \sum_{j=k+1}^d \mathbf{u}_j^T \Sigma \mathbf{u}_j$
 - Eigen vector: $\Sigma \mathbf{u} = \lambda \mathbf{u}$ eigen value
 - Minimizing reconstruction error equivalent to picking $(\mathbf{u}_{k+1}, \dots, \mathbf{u}_d)$ to be eigen vectors with smallest eigen values
- Memory lane: $\Sigma \mathbf{u} = \lambda \mathbf{u}$ eigen value
- $\text{error}_k = \sum_{j=k+1}^d \mathbf{u}_j^T \Sigma \mathbf{u}_j$ is an eigen vector
- $\mathbf{u}_j^T \Sigma \mathbf{u}_j = \mathbf{u}_j^T \lambda_j \mathbf{u}_j$ (orthogonal)
- $= \lambda_j \mathbf{u}_j^T \mathbf{u}_j$
- $= \lambda_j$

$\min_{\mathbf{u}_1, \dots, \mathbf{u}_k} \text{error}_k \equiv$ throwing out data... and
with smallest eigen values of Σ

\equiv keeping $\mathbf{u}_1, \dots, \mathbf{u}_k$ with highest eigen values

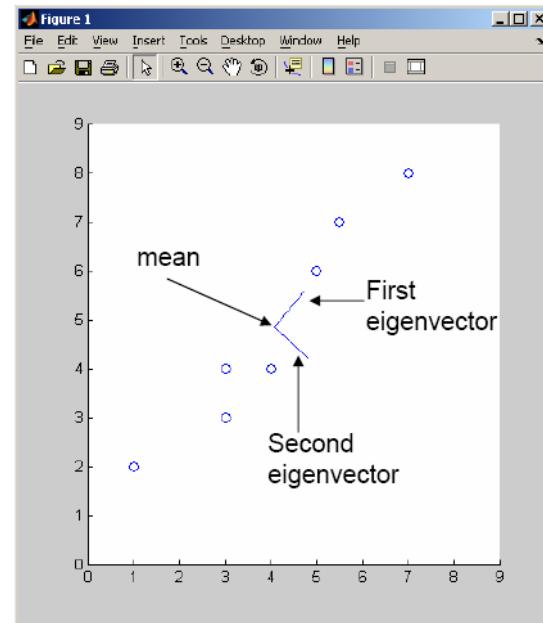
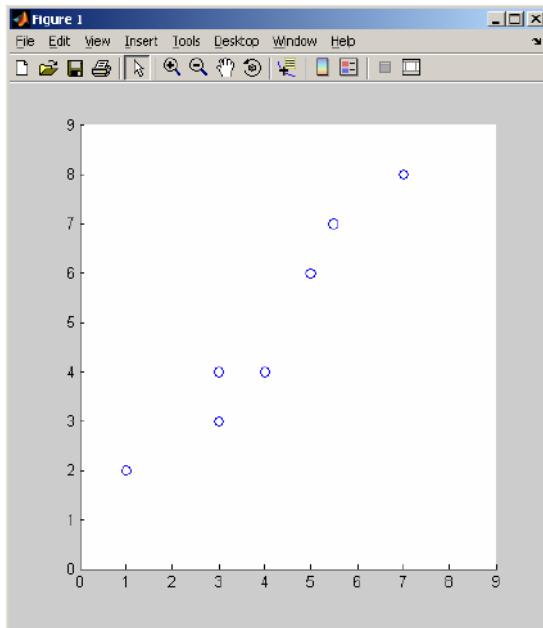
Basic PCA algorithm

- Start from N by d data matrix X
- Recenter: subtract mean from each row of X
 - $X_c \leftarrow X - \bar{X}$
- Compute covariance matrix:
 - $\Sigma \leftarrow 1/N X_c^T X_c$
- Find eigen vectors and values of Σ
- Principal components: k eigen vectors with highest eigen values

$$X_c = N \begin{pmatrix} d \\ \vdots \\ 1 \end{pmatrix} \begin{pmatrix} x_1 & \dots & x_d \end{pmatrix}^T$$

PCA example

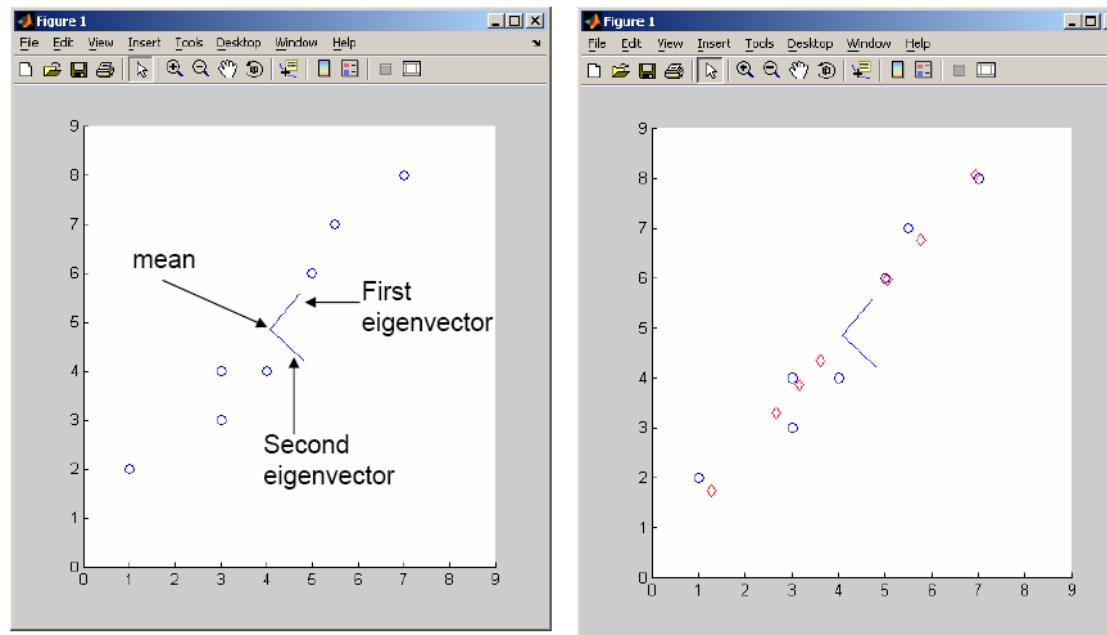
$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$



PCA example – reconstruction

$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$

only used first principal component



Eigenfaces [Turk, Pentland '91]

- Input



Glasses

- Principal components:



Eigenfaces reconstruction

- Each image corresponds to adding 8 principal components:



Scaling up

- Covariance matrix can be really big!
 - Σ is d by d
 - Say, only 10000 features
 - finding eigenvectors is very slow...
- Use singular value decomposition (SVD)
 - finds k eigenvectors
 - great implementations available, e.g., python, R, Matlab svd

SVD

- Write $X = W S V^T$
 - $X \leftarrow$ data matrix, one row per datapoint
 - $W \leftarrow$ weight matrix, one row per datapoint – coordinate of x^i in eigenspace
 - $S \leftarrow$ singular value matrix, diagonal matrix
 - in our setting each entry is eigenvalue λ_j
 - $V^T \leftarrow$ singular vector matrix
 - in our setting each row is eigenvector v_j

PCA using SVD algorithm

- Start from m by n data matrix X
- **Recenter:** subtract mean from each row of X
 - $X_c \leftarrow X - \bar{X}$
- Call SVD algorithm on X_c – ask for k singular vectors
- **Principal components:** k singular vectors with highest singular values (rows of V^T)
 - **Coefficients** become:

What you need to know

- Dimensionality reduction
 - why and when it's important
- Simple feature selection
- Principal component analysis
 - minimizing reconstruction error
 - relationship to covariance matrix and eigenvectors
 - using SVD